



構造/テキストWebデータを対象とした Pay-as-you-goスタイルの問合せ構築支援

著者	安永 ゆい
内容記述	筑波大学修士（情報学）学位論文・平成25年3月25日授与（30971号）
発行年	2013
URL	http://hdl.handle.net/2241/121397

構造/テキスト Web データを対象とした
Pay-as-you-go スタイルの問合せ構築支援

筑波大学
図書館情報メディア研究科
2013 年 3 月
安永 ゆい

目次

第 1 章	はじめに	1
第 2 章	関連研究	4
2.1	構造化問合せを対象とした問合せ記述支援に関する研究	4
2.2	問合せ記述支援に関する様々なアプローチ	5
第 3 章	Gradation 問合せ言語	6
3.1	特徴	6
3.2	問合せ対象データ	7
3.3	問合せ処理モデル	8
3.4	Gradation の問合せ構成要素と問合せ例	8
第 4 章	提案手法 MorphingAssist	10
4.1	取り組んだ問題	10
4.1.1	表示するヒント	11
4.1.2	ヒントの計算方法	11
4.1.3	ヒントの提示方法	11
4.2	ヒントの生成	12
4.2.1	【Phase1】使われる可能性のある全ての問合せ構成要素を求める	13
4.2.2	【Phase2】問合せ構成要素を使って問合せ改訂案を生成する	14
第 5 章	評価実験	16
5.1	データと被験者	16
5.2	方法	16
5.2.1	情報要求と正解問合せの作成	17
5.2.2	正解問合せの構成要素の作成	18
5.2.3	被験者による問合せ記述	18
5.3	実験結果 1：表示したヒントの質	18

5.4	実験結果 2 : MorphingAssist の効果	20
5.4.1	確定問合せの検索結果の再現率と精度	20
5.4.2	ユーザが書き直した問合せがより精確なものになっているか	20
第 6 章	終わりに	24
参考文献		25
付録 A	実験で用いた情報要求と正解問合せ	29
付録 B	スコア $ds_{cq}(q)$ の求め方	30
B.1	構成要素の対応付け関数	30
B.2	ペナルティの値の決定	31

第 1 章

はじめに

近年、テキストデータとそれと並存する構造データから構成された Web データが広く普及しつつある。このような Web データの例として、Wikipedia のページとそれと並存する DBpedia [1] のデータの組が挙げられる。DBpedia は Wikipedia の各記事に関する様々な情報を RDF 形式の構造データとして提供している。もう 1 つの例としては、ACM Digital Library [2] のページと、DBLP RKB Explorer [3] のデータとの組が挙げられる。ACM Digital Library はコンピュータサイエンス分野の論文の書誌情報をテキストページで提供しており、DBLP RKB Explorer は、コンピュータサイエンス分野の論文書誌データベースである DBLP のデータを RDF 形式で持っている。各書誌情報のページと RDF データ間の関係は論文に付与された DOI で表される。Semantic Web ビジョンにおける Linked Open Data の取り組み [4] 等と相まって、テキストデータとそれと並存する構造データから構成された Web データは今後ますます一般的に利用されと考えられる。

これまで、先述のような Web データに問合せを行う方法には、2 つの選択肢しか存在しなかった。すなわち、Web 検索などで使われるキーワード問合せを用いるか、もしくは SQL や SPARQL などの構造化問合せ言語に従って記述する構造化問合せを用いるか、のどちらかである。キーワード問合せは、カジュアルユーザに普及しているが、表現力が低く、Web データのセマンティクスを利用した問合せを記述することができない。一方、構造化問合せは、表現力が高く高度な問合せが可能であるが、言語の学習コストが大きく、また、問合せ対象のデータに関する知識も必要となるため、カジュアルユーザが習得・利用することは難しい。したがって、ユーザにとっては、テキストデータを対象としたキーワード問合せを行うか、構造データを対象とした構造化問合せを行うか、という二者択一の状況であった。

既存研究では、この問題に対する新しいアプローチとして、テキストページ (本論文では、テキストを含む Web ページなどをテキストページと呼ぶ) とグラフデータ (本論文では、主に RDF データを指す) の組からなる Web データに対する問合せ言語である Gradation 問合せ言語 (Gradation Query Language. 以下、Gradation) が提案されている [5][6]。Gradation の設計目標は、キーワード問合せと構造化問合せの溝を埋めることによって、高度な問合せ能力をカジュアルユーザにとって身近なものにすることである。Gradation は、キーワード問合せと構造化問合せとのシームレスな

融合を実現するために、キーワード問合せをベースとし、簡単な追加記述によって構造データに対する問合せ条件を表現する。**Gradation**においては、キーワードのみを使った問合せは通常のキーワード問合せのように動作し、構造データに対する問合せ条件のみを使った問合せは完全な構造化問合せとなり、キーワードと構造データに対する問合せ条件を混在させた問合せはテキストとグラフデータに対するハイブリッド問合せとなる。このような問合せ記述方法を採用することにより、**Gradation**は単純なキーワード問合せから高度な構造化問合せまでを広くカバーするとともに、“pay-as-you-go”スタイルの問合せ記述を目指す。すなわち、問合せ記述にかかるコストや問合せ条件の精確さを、ユーザが自身の要求やスキルに応じて選択し、利用することを可能にしようとしている。しかし、**Gradation**で構造データに対する問合せ条件を記述する（つまりハイブリッド問合せや完全な構造化問合せを記述する）ためには、検索対象データのクラス名や属性名などのスキーマレベルの情報が必要であるため、SQLやSPARQLといった他の構造化問合せ言語を使用する場合と同じだけの知識が求められる。よって、実際には**Gradation**だけでは“pay-as-you-go”スタイルの問合せ記述を実現することはできない。

本論文では、**Gradation**で構造データに対する問合せ条件を記述することを支援する手法 **MorphingAssist**を提案する。**MorphingAssist**は、**Gradation**の「キーワード問合せから構造化問合せまでを同一の言語でシームレスに記述できる」という特徴に着目した、より精確 (precise) な問合せの作成支援手法である。本支援手法の目標は、ユーザがキーワード問合せによる問合せから始め、**MorphingAssist**が提示するヒントを見ながら徐々に問合せを「必要な範囲で」「出来る範囲で」詳細化することを支援することである。本支援手法により、ユーザは完全な構造化問合せやハイブリッド問合せを容易に記述できる。ユーザは必ずしも完全な構造化問合せを作成する必要はなく、必要な範囲、出来る範囲で詳細化したハイブリッド問合せを記述することができる。この点において **MorphingAssist**は“pay-as-you-go”スタイルの問合せ構築を支援すると言える。

MorphingAssistの開発にあたっては、自明でない3つの問題に取り組む必要がある。(1) **表示するヒント**: どのようなヒントを提示すればよいのか。(2) **ヒントの計算方法**: どのようにしてそのヒントを求めるか。(3) **ヒントの提示方法**: どのように求めたヒントをユーザに見せるべきか。詳細は4章で説明するが、それぞれの問題について次のようなアプローチを採用する。(1)に関しては、一般的なキーワード検索の **Suggest** 機能と同様に、「もしかしたら、この(構造データに関する)問合せ構成要素を使いたかったのかもしれない」という問合せ改訂案を見せるというアプローチを採用する。(2)に関しては、ユーザが文書検索を意図して入力したキーワードをあえて構造データとマッチさせ、ユーザが使用する可能性のある問合せ構成要素を計算するというアプローチを採用する。(3)に関しては、ヒントの組み合わせ爆発を避けつつ、ユーザが問合せの書き直しをしやすいと考えられる手法を採用する。実験の結果、提案アプローチは容易に実現可能であるにも関わらず効果的であることがわかった。

本論文の貢献は次の通りである。

(1) **二択でない構造化問合せの作成支援**. 「完全な構造化問合せを出力できるか、そもそも問合せが

出来ないかの二択」であった既存の SQL 言語等の構造化問合せ作成支援ツールと異なり，Gradation 問合せ言語を利用する事によって，キーワード問合せから徐々に「必要な範囲で」「出来る範囲で」問合せの詳細化を支援するという，新規性の高いアプローチを提案する．

(2) 実データによる評価実験．Wikipedia/DBpedia を用いて被験者を用いた評価実験を行った結果を示す．実験では，Gradation 言語と MorphingAssist を利用する事により，SQL 言語等の構造化問合せを書けない，もしくは書いたことがない被験者の多くが精確な問合せを記述することができた．

本論文の構成は次の通りである．2 章では，関連研究について述べる．3 章では，キーワード問合せと構造化問合せをシームレスに融合した問合せ言語 Gradation の概要を説明する．4 章で，Gradation を用いた構造問合せ作成支援 MorphingAssist を提案する．5 章では評価実験について述べる．6 章では本論文全体のまとめについて述べる．

第 2 章

関連研究

本研究は、構造化問合せ言語で問合せを記述することを支援する研究である。本研究の関連研究は大きく次の 2 つに分けられる。(1) 構造化問合せを対象とした問合せ記述支援に関する研究。(2) 問合せ記述支援に関する様々なアプローチ。

2.1 構造化問合せを対象とした問合せ記述支援に関する研究

ユーザが構造化問合せ言語で問合せを記述することを支援するシステムは数多く存在し、大きく次の 2 つに分類される。(1) ユーザに構造化問合せを入力させ、よりユーザの意図に合致と思われる構造化問合せを構築していく支援方法。(2) ユーザにキーワードや自然言語文を入力させ、構造化問合せを構築する支援方法。

(1) の例としては、ユーザの書いた SQL 問合せを元に別の SQL 問合せを推薦する研究 [8] が存在する。[8] では、ユーザは初めに構造化問合せを記述する必要がある。これに対して、提案手法は、構造化問合せが記述できないユーザが構造化問合せを記述することを支援することを目的としているので、ユーザは初めにキーワード問合せを入力するだけでよい。

(2) の例としては IQP [7] がある。IQP は最初にユーザに問合せを表すキーワードを入力させ、次にユーザに質問を行い、最終的にユーザの意図を反映した構造化問合せを出力する。本研究との違いは次の 2 点である。第 1 に、IQP は最終的な構造化問合せが出力するまで、ユーザが問合せを実行することが出来ない。これに対して、Gradation を利用する本提案手法では、完全な構造化問合せでないキーワード問合せやハイブリッド問合せであっても実行することができる。これにより、ユーザが自身の要求やスキルに応じた程度の問合せを記述し実行可能である。第 2 に、IQP では、最初に入力されたキーワードが、最終的に出力される構造化問合せに大きく影響する。したがって、最初に適切なキーワードを入力できなければ、生成される問合せはユーザの意図を反映しなくなる可能性が高い。しかし、我々の経験から言うと、ユーザが最初から適切なキーワードを選ぶことは困難であることが多い。これに対して、Gradation を利用すれば、最初にキーワード問合せを実行したときから検索結果がその場で表示されるので、記述したキーワード問合せが不適切だと気付けば、

その場で簡単にキーワードを変更することができる。

2.2 問合せ記述支援に関する様々なアプローチ

問合せ記述支援に関するアプローチは、大きく次の2つに分類される。(1) ユーザが入力した問合せを元に、新しい問合せを推薦する方法。(2) ユーザとインタラクティブなやりとりを行い、ユーザの問合せ構築を支援する方法。

(1) の手法は Web 検索における **suggest** 機能をはじめとし、様々なシステム・研究で用いられている。既存研究の多くは、ユーザがキーワード問合せを入力した際にキーワード問合せを推薦するもの、もしくはユーザが構造化問合せを入力した際に構造化問合せを推薦するものである。本研究も問合せを推薦するが、本研究はキーワード問合せを入力とし、構造データに関する条件を推薦する点で、多くの既存研究と異なる。

(2) の手法をとっている研究の例としては、2.1 節で述べた **IQP** が挙げられる。**IQP** は各インタラクションで質問を提示し、全てのインタラクションが終わって初めて構造化問合せが構築できる、という手法である。本研究もユーザとのインタラクティブなやりとりによりユーザの問合せ構築を支援するものであるが、本研究は毎インタラクションで、直前にユーザが入力した問合せをより詳細化した問合せを推薦するので、ユーザは自分の要求やスキルに応じた問合せが構築できれば、それ以上インタラクションを続ける必要はない。

第 3 章

Gradation 問合せ言語

本章では Gradation の概要を説明する．まず Gradation の特徴を説明する．次に Gradation の問合せ対象データと Gradation の問合せ処理モデルを説明する．最後に Gradation の問合せ構成要素を問合せ例を用いて説明する．

3.1 特徴

Gradation は、テキストページと構造データから構成された Web データを検索対象としたハイブリッド問合せ言語である．Gradation は、キーワード問合せをベースとし、簡単な追加記述によって構造データに対する問合せ条件を表現する言語であり、キーワード問合せと構造化問合せとのシームレスな融合を実現している言語である．

Gradation は、キーワード問合せのみを利用するようなユーザであっても、キーワードに加えて簡単な検索オプションをしばしば利用することに着目して設計されている．例えば、Web 検索エンジンのユーザは、指定したサイトのみを検索対象とするオプション (`site:ac.jp` など) を記述して検索することがある．実際、[9] によると、25% のユーザが検索オプションを使用したことがあると言われている．

Gradation でも、キーワードに加えて検索オプションのような簡単な追加記述を利用することができる．この簡単な追加記述によって構造化データに関する問合せ条件を表現する．例えば、グラフデータに各人物の年齢 (`age`) を表すデータが含まれる場合、40 歳以上の Tom という人物のテキストページを取得する問合せは、追加記述 `age>=40` を用いて `Tom age>=40` と記述する．

Gradation は単なるキーワード問合せや構造化問合せが記述できるだけでなく、ハイブリッドな問合せを記述することができる．ハイブリッドな問合せというのは、キーワードと構造化データに関する問合せ条件を混合して記述している問合せである．例えば、先に示した問合せ `Tom age>=40` はハイブリッドな問合せである．

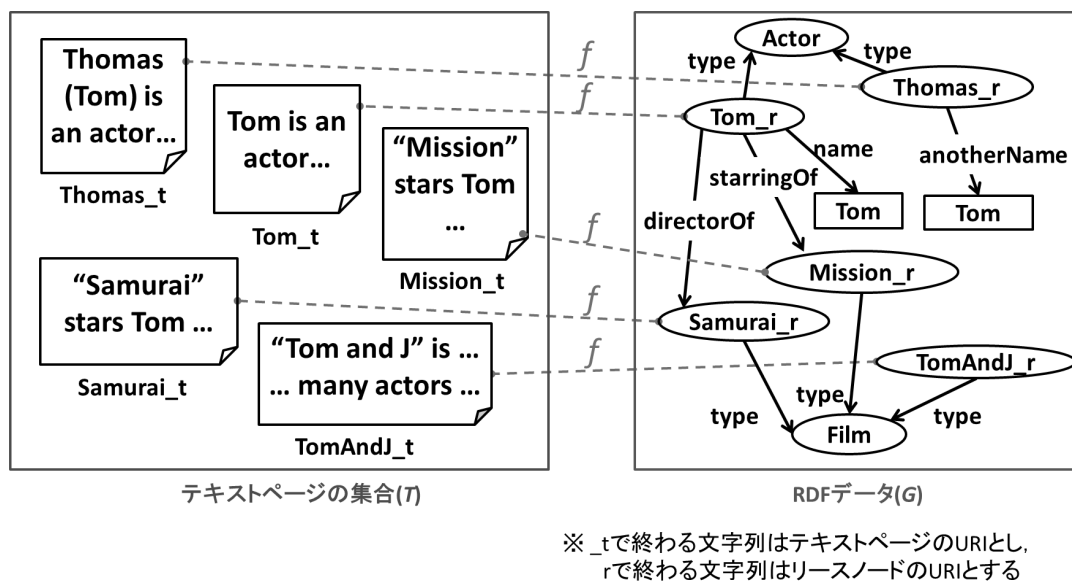
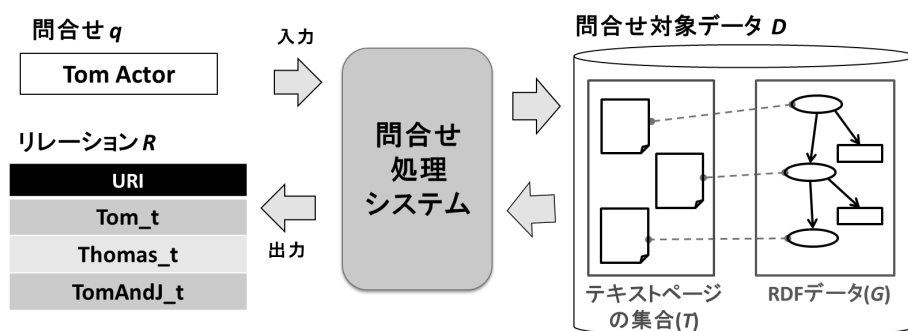
図 3.1 問合せ対象データ D の例

図 3.2 問合せ処理モデル

3.2 問合せ対象データ

Gradation の問合せ対象データ D は、トリプル $D = (T, G, f)$ で表現される。図 3.1 は問合せ対象データの一例である。 T はテキストページの集合、 G はテキストページの集合と並存するグラフデータ (RDF データ)、 f はテキストページと RDF 中のノードの対応関係を表す関数 $f: T \rightarrow G$ を表す。図 3.1 を含め、本論文では、 $_t$ で終わる文字列はテキストページの URI とし、 $_r$ で終わる文字列はリソースノードの URI とする。

3.3 問合せ処理モデル

Gradation の問合せ処理モデルを図 3.2 に示す。入力は Gradation で記述した問合せ q (図 3.2 左上), 問合せ対象データは 3.2 節で説明した問合せ対象データ D (図 3.2 右), 出力は, 特別な指定がなければ, q にマッチするテキストページの URI の集合を表すリレーション R (図 3.2 左下) である。特別な指定については 3.4 節 (C4) で述べる。

3.4 Gradation の問合せ構成要素と問合せ例

本節では Gradation の問合せ構成要素 (C1) から (C5) とそれを使った問合せ例を説明する。下記問合せ例では, 図 3.1 を問合せ対象データとする。

(C1) keyword: あるキーワードを含むテキストページを探すための構成要素。例えば, Tom と Actor という文字列を含むテキストページを取得する問合せは次のようになる。

Tom Actor

この問合せの検索結果は図 3.3(C1) の通りである。

(C2) selection: 属性を指定する構成要素。“ $attr \theta value$ ” もしくは “ $attr_1 \theta attr_2$ ” のように記述する。前者は, RDF データ中でトリプル $(v, attr, value)$ が存在し, 後者は, RDF データ中でトリプル $(v, attr_1, value_1)$, $(v, attr_2, value_2)$ (ただし $value_1 = value_2$) が存在することを表す。 θ には, 等号・不等号の記号が使える。例えば, 名前 (name) が Tom である俳優のテキストページを取得する問合せは次のようになる。

name==Tom type==Actor

この問合せの検索結果は図 3.3(C2) の通りである。

(C3) path traversal: 関係を指定する構成要素。 $q_1.p.q_2$ のように記述する。 q_1 と q_2 には任意の構成要素を記述することができる。 q_1 が示す RDF 中のノード v_1 と, q_2 が示す RDF 中のノード v_2 の間にエッジ p が存在することを表す。例えば, 名前が Tom である俳優と, 彼が主演を務めた (starringOf) 映画のテキストページの組を取得する問合せは次の通りである。

(name==Tom type==Actor).starringOf.type==Film

この問合せの検索結果は図 3.3(C3) の通りである。本構成要素を使用すると, 検索結果のリレーションには, p の関係にあるテキストページの組が表示される。

(C4) projection: 検索結果に表示する属性名を指定する構成要素。 $q [attr_1, \dots, attr_n]$ のように記述する。 q には任意の構成要素を記述することができる。 $attr_i$ には RDF データ中のエッジ名の他, テキストページの URI を表す t_uri と RDF データのリソースノードを表す r_uri が記述できる。例えば, 俳優を取得し, その俳優のテキストページ (t_uri) と名前 (name) を表示する問合せは次の通

(C1):			(C4):	(C5):
uri	(C2):	(C3):	t_uri	uri
Thomas_t	uri	uri(Actor)	Tom_t	Thomas_t
Tom_t	Tom_t	uri(Film)	Thomas_t	Tom_t
TomAndJ_t		Tom_t		
		Samurai_t		

図 3.3 構成要素 C1 から C5 それぞれの問合せ例の実行結果

りである。

```
(type==Actor)[t_uri, name]
```

この問合せの検索結果は図 3.3(C4) の通りである。

(C5) and, or, not, *, (): 任意の構成要素に対して、and, or, not, *(直積), グループングを行うことができる。例えば、俳優 (Actor) である人物のテキストページで、かつ Tom という文字列を含むテキストページを取得する問合せは次の通りである (ただし構成要素 “and” は省略可能である)。この問合せは、キーワード問合せの条件と構造化問合せの条件が混在した問合せである。

```
Tom and type==Actor
```

この問合せの検索結果は図 3.3(C5) の通りである。

以上の問合せ例のように、Gradation を用いると、キーワードと構造データに関する問合せ条件を混在させ、簡易なレベルの問合せから高度な問合せまでが表現可能である。

第 4 章

提案手法 MorphingAssist

3 章で述べたように、**Gradation** を利用すれば、キーワード問合せから構造化問合せまでを同一の言語でシームレスに記述できる。しかし、構造化データに関する問合せ条件 (すなわち構成要素 **C2** から **C4**) を記述するためには、対象データのクラス名や属性名などのスキーマレベルの情報が必要であるため、**SQL** や **SPARQL** といった他の構造化問合せ言語を使用する場合と同じだけの知識が求められる。

本章では、**Gradation** での精確な問合せ作成を支援する手法 **MorphingAssist** を提案する。これは、**Gradation** の「キーワード問合せから構造化問合せまでを同一の言語でシームレスに記述できる」という特徴に着目した、精確な問合せの作成支援手法である。**MorphingAssist** のユーザはキーワード問合せから始め、その結果を見ながらキーワードを変更したり、**MorphingAssist** の支援によって徐々に問合せを「必要な範囲で」「出来る範囲で」詳細化したりする。例えば、「俳優のトムという人物を全員知りたい」と思ったユーザが最初に問合せ (QI) “Tom Actor” を記述したとする。もし **MorphingAssist** が適切な支援を行えば、ユーザは徐々に問合せをより精確な問合せに書き直して、(QII) “Tom type==Actor” や (QIII) “name==Tom type==Actor” を記述できる。**MorphingAssist** はこのような問合せの変更を支援するためのヒントを提示する。

4.1 取り組んだ問題

MorphingAssist の開発にあたって、自明でない次の 3 つの問題に取り組んだ。(1) **表示するヒント**: どのようなヒントを提示すればよいのか。(2) **ヒントの計算方法**: どのようにしてそのヒントを求めるか。(3) **ヒントの提示方法**: どのように求めたヒントをユーザに見せるべきか。本節ではそれぞれの問題へのアプローチを説明する。

4.1.1 表示するヒント

キーワード検索システムの問合せ推薦機能の一般的なアプローチは、「もしかしたら、このキーワードを使いたかったのかもしれない」と考えられる新しいキーワードを次の問合せのために推薦する、というアプローチである。本研究では、このアプローチを本研究の文脈に適用し、「もしかしたら、この(構造データに関する)問合せ要素を使いたかったのかもしれない」と考えられる問合せ要素を次の問合せのために推薦するというアプローチを採用する。本アプローチでは、キーワード問合せを書いたユーザに対して構造データに関する問合せ構成要素を推薦する。これにより、キーワード問合せを書いたユーザは推薦された構成要素を使うと元の問合せをより精確な問合せに書き直すことができる。例えば、(QI)を書いたユーザに対して問合せの構成要素“name==Tom”や“type==Actor”を推薦する。この構成要素を使うと、ユーザは(QI)を書き直して(QII)や(QIII)を書くことができる。

4.1.2 ヒントの計算方法

ヒントを求めるために、本研究ではユーザが書いた元のキーワード(テキストデータとマッチすることを意図して書かれたキーワード)をあえてそれと並存する構造データにマッチさせ、ユーザが使いたかった可能性のある問合せ構成要素を計算するというアプローチを採用する。この手法は単純であり、同時に、効果的であることが実験(5章参照)によって示されている。

4.1.3 ヒントの提示方法

ヒントの単純な提示方法の1つは、ユーザが使いたかった可能性のある問合せの構成要素を並べて提示することである。しかし、これではユーザは提示された構成要素をどのように使えばいいか思いつくことができない。別の単純な方法としては、元の問合せと求めた構成要素を使って作成可能なより詳細な問合せへの改訂案(以下、**問合せ改訂案**)を**全て**提示することである。しかし、これでは提示される問合せ改訂案の数は膨大になり、容易に組合せ爆発を起こしてしまう。

そこで、本研究ではヒントの組み合わせ爆発を避けつつ、ユーザが問合せの書き直しをしやすいと考えられる提示方法を採用する。その方法とは、元の問合せと1カ所だけ異なる問合せ改訂案を全て列挙することである。こうすることにより、ヒントの数は求めた構成要素の数と同じになり、かつユーザは提示されたヒントを組み合わせることで2つ以上の構成要素を合わせて使った新たな問合せを簡単に思いつくことができる。例えば、元の問合せが“Actor Film Star”であり、求めた構成要素が“type==Actor”, “type==Film”, “[star]”, “.star.”の4つだったとする。この時、MorphingAssistは次の4つのヒントを提示する。“type==Actor Film Star”, “Actor

URL	Hints
Tom_t	<ul style="list-style-type: none"> • name==Tom Actor • Tom type==Actor
Thomas_t	<ul style="list-style-type: none"> • anotherName==Tom Actor • Tom type==Actor
TomAndJ_t	-----

図 4.1 検索結果とヒントの例

type==Film Star”, “(Actor Film)[star]”, “(Actor Film).star.(*)^{*1}”.

MorphingAssist では、直前に投げた問合せの検索結果の 1 タプル毎にヒントが作成される。例えば、ユーザが問合せ (QI) を図 3.1 のデータに対して投げると、図 4.1 のように検索結果とヒントが表示される。著者の作成するユーザインタフェース (図 4.2) では、ユーザは各タプルの隣にある「ヒントを見る」ボタンをクリックするとそのタプルに関連するヒントが表示されるようにしている。これは、ユーザの関心のない検索結果 (タプル) に関連するヒントはユーザに必要なと考えたためである。図 4.1 の検索結果とヒントを見たユーザは、1 行目のヒントを見れば問合せ (QII) が書け、1 行目と 2 行目のヒントを組み合わせれば問合せ (QIII) が書ける。なぜなら、1 行目と 2 行目のヒントはそれぞれ、元の間合せの “Tom” は “name==Tom” に置き換えることができ、“Actor” は “type==Actor” に置き換えることができることを示しているからである。

表示数を減らす仕組み。 さらに、適切そうなヒントを選び出すために、この提案手法を次の単純な手法と組み合わせる。まず、ユーザに、情報要求にマッチしているタプルを選ばせる。次に、それらのタプルと関連するヒントの積集合を計算し、ユーザに見せる。こうすることで、ユーザの情報要求を問合せで表現する際に重要なヒントがこの中に含まれることになると考えられる。図 4.3 に、実験で用いたインタフェースを示す。図 4.3 では、ユーザに、情報要求にマッチしているタプルの左のチェックボックス (図 4.3(1)) にチェックを入れさせ、ページ下部の「共通するヒントを見る」(図 4.3(2)) をクリックさせる。システムはその情報を元に、ユーザがチェックを入れたタプルと関連するヒントの積集合を表示 (図 4.3(3)) する。

4.2 ヒントの生成

本節では MorphingAssist が生成するヒントの生成手順を説明する。ヒントの生成は次の 2 段階からなる。【Phase1】使われる可能性のある全ての問合せ構成要素を求める。【Phase2】問合せ構成要素を使って問合せ改訂案を生成する。

^{*1} Gradation ではキーワードのワイルドカードとして “*” を使うことができる。



図 4.2 著者の作成したヒントの表示方法のインターフェース

4.2.1 【Phase1】使われる可能性のある全ての問合せ構成要素を求める

Gradation 問合せ処理システムによって求めた q_i の検索結果を R とする。【Phase1】では、検索結果 1 件 r ($r \in R$) 毎に、問合せ改訂案 q_{i+1} で使われる可能性のある構成要素を求める。すなわち、図 4.1 で言うと、各タプル毎に構成要素は求められる。具体的には、【Phase1】は次の 2 ステップからなる。

Step 1: 検索結果 1 件 r ($r \in R$) に含まれるテキストページ t の URI (検索結果のうち属性 t_uri) の集合を U とする。このとき、 $V = \{f(t) | t \in U\}$ となるような URI の集合 V を求める。ここで、関数 $f(t)$ は t と対応する RDF ノードを返す関数である (3 章参照)。

Step 2: V 中の RDF ノードから距離 h 以内にある RDF ノードの集合を V' とする。また、 V' を求める際にトラバースしたエッジのエッジラベルの集合を L とする。このとき、次のルールに基づいて構成要素を生成する。

Rule1: 問合せ q_i 中のキーワード 1 つが V' 中のノードのノードラベルとマッチしたら、selection 構成要素 (C2) “ $attr = value$ ” を生成する。ここで、 $value$ はマッチしたノードラベルの値であり、 $attr$ はそのノードから出ているエッジのエッジラベルである。もしマッチしたノードに対してそのノードから出ているエッジが複数存在したら、それぞれのエッジに対して 1



図 4.3 著者の作成した表示数を減らす仕組みのインターフェース

つ構成要素を生成する．加えて，もし *value* が数字であった場合は，不等号記号を使った *selection* 構成要素 (“*attr* <= *value*” など) も生成する．

Rule2: 問合せ q_i 中のキーワード 1 つが L 中のエッジラベルとマッチしたら，*projection* の構成要素 (C4) “[*l*]” と *path-traversal* の構成要素 (C3) “*l*.” を生成する．ここで，*l* はマッチしたラベルである．

例えば，図 4.1 の 1 タプル目の検索結果 (Tom.t) に対応する構成要素は次のようにして求める． q_i は “Tom Actor” で，ヒントを求めたい検索結果 (1 タプル目) は *uri* 属性に Tom.t を持っている．このとき，図 3.1 では $f(\text{Tom.t}) = \text{Tom.r}$ が成り立つ．図 3.1 では，Tom.r の隣接ノードでキーワード “Tom” とマッチするノード (“name” エッジで繋がっているノード) とキーワード “Actor” とマッチするノード (“type” エッジで繋がっているノード) がある (ここで，“Tom” と “Actor” は問合せ q_i に含まれているキーワードである)．従って，図 4.1 中の 1 タプル目の検索結果 (Tom.t) に対応する構成要素は “name==Tom” と “type==Actor” の 2 つとなる．

4.2.2 【Phase2】問合せ構成要素を使って問合せ改訂案を生成する

ここでは，(1) で生成した問合せ構成要素の 1 つと q_i 中のキーワード 1 つを置き換えることでヒント (問合せ改訂案) を生成する．例えば， q_i が Tom Actor で 1 件目の検索結果に対して生成され

た問合せ構成要素が “name==Tom” と “type==Actor” の 2 つだった場合, (2) で生成するヒントは “name==Tom Actor” と “Tom type==Actor” の 2 つである.

第 5 章

評価実験

本章では MorphingAssist の評価実験について述べる。本実験では次の 2 つの観点から評価を行った。(1) 生成したヒントの質。(2) MorphingAssist の効果。

5.1 データと被験者

本実験で使用したデータは、Wikipedia [10] と DBpedia Japanese [11] のダンプデータである。

本実験では、2 つのグループの被験者に別々の作業をしてもらい実験を行った。**グループ A** : 実験で使用する情報要求を作成する学生 (7 名)。**グループ B** : MorphingAssist システムを使って Gradation 問合せを記述する学生 (10 名)。さらにグループ B は、グループ B1 とグループ B2 に分けることができる。グループ B1 は DB 言語 (例えば SQL や XQuery, SPARQL など) をある程度使うことができる学生 (5 人) であり、グループ B2 は DB 言語を 1 つも知らない、もしくはかつて習ったことがあるが忘れてしまった学生 (5 人) である。グループ B 全員は Gradation の文法に関する知識がなかったため、実験前に文法の説明を行った。

5.2 方法

本実験は次の手順で行った。

- (1) **情報要求と正解問合せの作成**: まず、グループ A の 7 人の被験者に情報要求 (Wikipedia や DBpedia に載っている情報に対する情報要求) を考えてもらった。この情報要求の集合を N とする。次に、著者がそれぞれの情報要求 $n_i \in N$ を満たすような正解 Gradation 問合せ cq_i を作成した。正解 Gradation 問合せの集合を CQ とする。
- (2) **正解問合せの構成要素を求める**: それぞれの正解 Gradation 問合せ $cq_i \in CQ$ ごとに、その問合せの構成要素の集合 C_i を求めた。例えば、 $cq_i = \text{“name==tom type==actor”}$ のとき、 $C_i = \{\text{name==tom, type==actor}\}$ である。
- (3) **情報要求から Gradation 問合せを記述**: グループ B の被験者に情報要求を提示し、MorphingAs-

要件	情報要求の例
物事の「属性」を使う	「芥川賞を受賞した」作家を全員知りたい.
複数の物事間の「関係」を使う	お笑い芸人と東京都に本社がある事務所の組み合わせ（「その芸人はその事務所に所属している」）全てが知りたい.
物事の「属性」を表示する	日本にある全ての動物園と「その開園年」が知りたい.

図 5.1 情報要求作成に係る 3 つの要件と情報要求の例

sist システムを使用して Gradation 問合せを記述してもらう．ここで、問合せ $q_{i,j,k}$ は、情報要求 $n_i \in N$ を見た被験者 s_j が書いた k 番目の問合せを表すこととする． $Q_{i,j}$ は、 $q_{i,j,k}$ を k を 1 から並べたリストとする．以降、 $Q_{i,j}$ を、問合せ記述プロセスと呼ぶ． $q_{i,j,k}$ は、 $Q_{i,j}$ の k 番目の問合せとなる．また、問合せ $q_{i,j,k}$ の結果と共に被験者 j に提示したヒントの集合を $H_{i,j,k}$ で表現する．

(4) ヒントの質と MorphingAssist の効果から MorphingAssist を評価: 4 章に示したように、MorphingAssist は、ユーザに再検索で使われる可能性のある問合せ構成要素を含む問合せをヒントとして提示する．本実験では、次の 2 つの観点で MorphingAssist を評価する．(1) 生成したヒントの質: ユーザに提示したヒントの表示数 $|H_{i,j,k}|$ 、ヒントの再現率 (正解問合せの構成要素 C_i をどれだけ再現できたか)．(2) MorphingAssist の効果: 被験者が記述した問合せの実行結果の再現率と精度、被験者が記述した問合せと正解問合せの比較．

5.2.1 節から 5.2.3 節で、手順 (1) から (3) の詳細を説明する．

5.2.1 情報要求と正解問合せの作成

まず、グループ A の被験者それぞれに情報要求 (Wikipedia や DBpedia に載っている情報に対する情報要求) を 3 つ自然言語で記述してもらった．このとき、各情報要求が図 5.1 に書いてある要件 3 つのうち少なくとも 1 つを満たし、「○○なものを全部知りたい」という形の情報要求になるように指示をした．これは、Gradation における構造化問合せ構成要素のうち、基本的な構成要素が次の 3 つであるためである．(1) 属性を指定する構成要素 (C2)、(2) 関係を指定する構成要素 (C3)、(3) 検索結果に表示する属性名を指定する構成要素 (C4)．情報要求を作成してもらうにあたって、被験者には情報要求のサンプルとして図 5.1 の情報要求の例を示した．

次に、集めた情報要求の中から、DBpedia 中に情報がないために調べることができない情報要求を排除した．その結果、10 個の情報要求が得られた (つまり、 $|N| = 10$)．

最後に、著者がそれぞれの情報要求 $n_i \in N$ を満たすような Gradation 正解問合せ $cq_i \in CQ$ を記述した．各 cq_i の実行結果を、情報要求 n_i を記述した被験者に見せ、その結果が n_i を満たしているかどうか確認してもらった．満たさない、と判断された場合は、満たすと判断されるまで cq_i を書き直した (本実験では、満たさない、と判断されることはなかった)．付録 A に、本実験で使った情報要求とそれを満たす正解 Gradation 問合せを示す．

5.2.2 正解問合せの構成要素の作成

まず、各正解問合せ $cq_i \in CQ$ から問合せの構成要素を抽出して、 C_i ($C_i = \text{component}(cq_i)$) を作成する。次に、 C_i から、単純な問合せ構成要素を削除する。単純な問合せ構成要素とは、RDF データ中のデータ (エッジラベル、ノードラベル) を含まないような構成要素である。具体的には、(C5) の “and”, “or”, “not”, “*” である。これらを削除した理由は、MorphingAssist は、RDF データ中のデータに関する情報をヒントとして提示するシステムであり、これらの構成要素が MorphingAssist によって提示されないからである。例えば、 cq_i が “(name==Tom or name==Thomas)” であったとき、 C_i は {name==Tom, name==Thomas} となる。

5.2.3 被験者による問合せ記述

それぞれの情報要求をグループ B の被験者 3 人に見せ、MorphingAssist システムを使って Gradation 問合せを記述してもらった。各被験者 s_j は 3 つの情報要求に対して Gradation 問合せを記述した。各被験者は、自分が書いた問合せ $q_{i,j,k}$ が情報要求 n_i を満たすと判断するまで問合せを書き直すことができる。被験者が自分で情報要求を満たすと判断した問合せを**確定問合せ**と呼び、 $dq_{i,j}$ と表す。被験者が 15 分以内に確定問合せを作成することができなかった場合は、確定問合せはなしとする。全ての被験者について最後の問合せ $q_{i,j,|Q_{i,j}|}$ は存在し、確定問合せが存在する場合は、 $dq_{i,j} = q_{i,j,|Q_{i,j}|}$ となる。

5.3 実験結果 1：表示したヒントの質

図 5.2 中の 1 つの点は問合せ記述プロセス 1 つ $Q_{i,j}$ を表す。ただし、ユーザが全くヒントを見なかった 2 プロセスは除く。y 軸は $|\cup_k H_{i,j,k}|/|Q_{i,j}|$ 、つまり各問合せ記述プロセス内で表示したヒントの数 (被験者がヒントを表示しなかった場合は除く) の平均を表す。 \cup_k は bug-union を求める記号を表し、 $\cup_k H_{i,j,k}$ は多重集合である。x 軸は C_i に対する ヒントの再現率を表す。ヒントの再現率とは、問合せ記述プロセス中で正解構成要素がどれだけ表示されたかを表すもので、 $|C_i \cap \cup_k H_{i,j,k}|/|C_i|$ によって求められる。 \cap は bug-intersection を求める記号を表すが、 C_i は多重集合ではないので、 $C_i \cap \cup_k H_{i,j,k}$ も多重集合ではなくなる。ヒントの再現率の平均は 0.71 であり、ヒントの表示数の平均は 25.76 であった。本評価観点で精度を求めないのは、確定問合せでない問合せが書かれた意図は一般的に不明であり、特定の意図をもって書かれた正解問合せに対する精度を求めることは意味をなさないからである。図 5.2 の全 28 プロセス中、4 プロセスで、再現率は 0.4 以下であった。しかし、これは被験者が MorphingAssist システムを使わなかったことが原因であった。

本実験では、30 問合せプロセス中、10 プロセスで、積集合を用いてヒントの表示数を減らす仕組み (4.1.3 節参照) が使われた。図 5.4 に、この 10 プロセスにおける、ヒントの表示数を減らす仕組

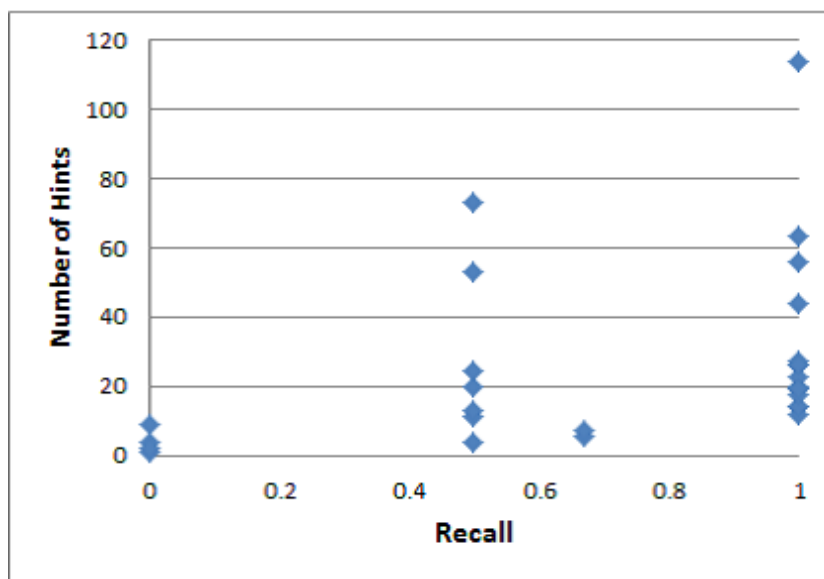


図 5.2 提案手法による構成要素の再現率とヒントの表示数

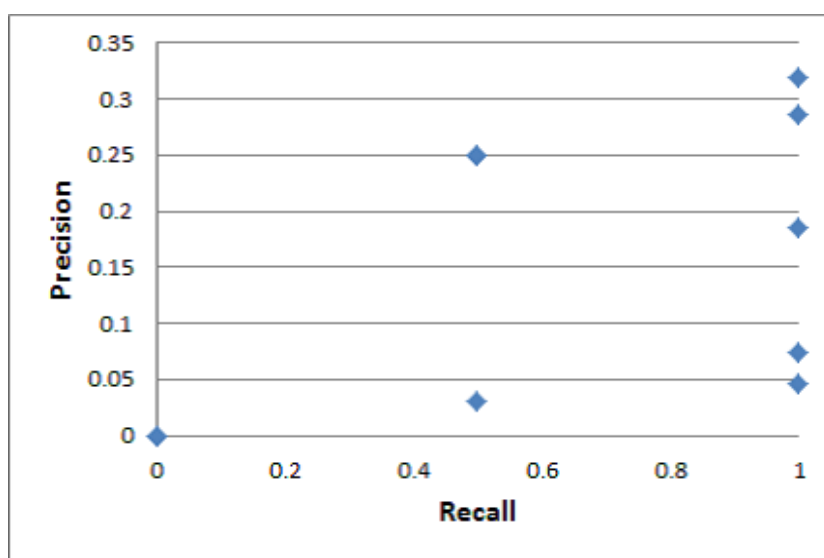


図 5.3 ヒントを減らす工夫をした時のヒントの再現率とヒントの精度

みを使う前と後の表示数の変化を示す。ヒントの表示数の平均は 43.22 から 9.71 に減少したにもかかわらず、再現率が下がったものは 1 つもなかった (再現率の平均は 0.65)。

被験者がヒントを減らす仕組みを使うということは、被験者の意図を表しているので、ここで精度を求めることは意味がある。図 5.3 は、ヒントを減らす仕組みを使用した後に表示されたヒントの再現率と精度である。ヒントの精度は、 $|C_i \cap \cup_k H_{i,j,k}| / |\cup_k H_{i,j,k}|$ によって求められる。ヒントの精度の平均は 0.16 であった。これは、表示されたヒント平均 9.71 個のうち 1.55 個は正解構成要

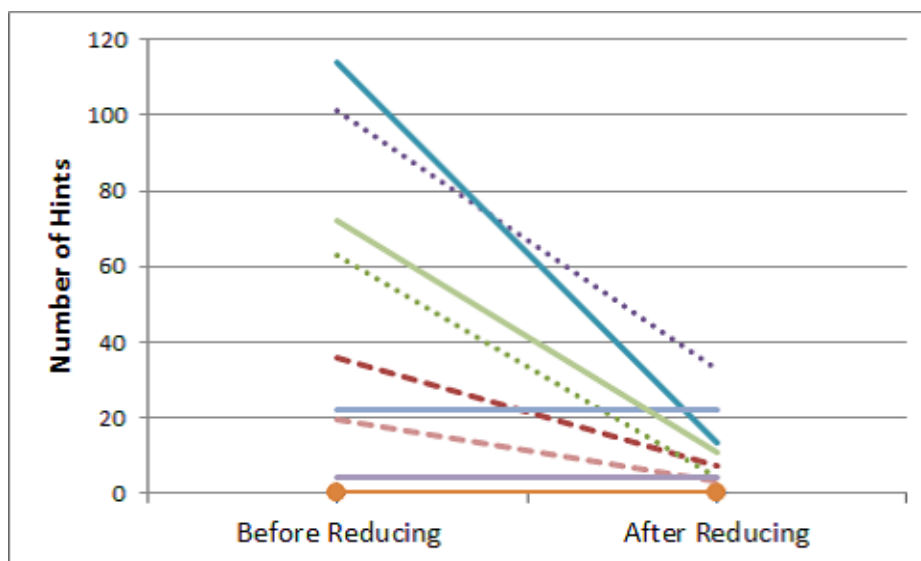


図 5.4 ヒントを減らす工夫をした時の表示されるヒント数の変化

素を含んでいることを示す。

以上の結果から、MorphingAssist は適切なヒントを表示することに関してよいはたらきをしていることが分かる。

5.4 実験結果 2 : MorphingAssist の効果

本節では、MorphingAssist の効果について次の 2 つの観点から評価する。(1) 確定問合せの検索結果の再現率と精度。(2)MorphingAssist の支援によりユーザが書き直した問合せがより精確なものになっているか。

5.4.1 確定問合せの検索結果の再現率と精度

図 5.5 に、正解問合せの検索結果に対する確定問合せの検索結果の再現率と精度を示す。再現率は $|Result(dq_{i,j}) \cap Result(cq_i)| / |Result(cq_i)|$ であり、精度は $|Result(dq_{i,j}) \cap Result(cq_i)| / |Result(dq_{i,j})|$ である。図 5.5 より多くの被験者は正しい問合せを記述することができたことが分かる。着目すべき点は、被験者のうちの半分は SQL を書くことができない人であったことである。

5.4.2 ユーザが書き直した問合せがより精確なものになっているか

次に、本節では確定問合せに向けて問合せがどのように変化したかを評価する。これは、前節で示した結果が本当に MorphingAssist によるものかどうかを確認するためである。

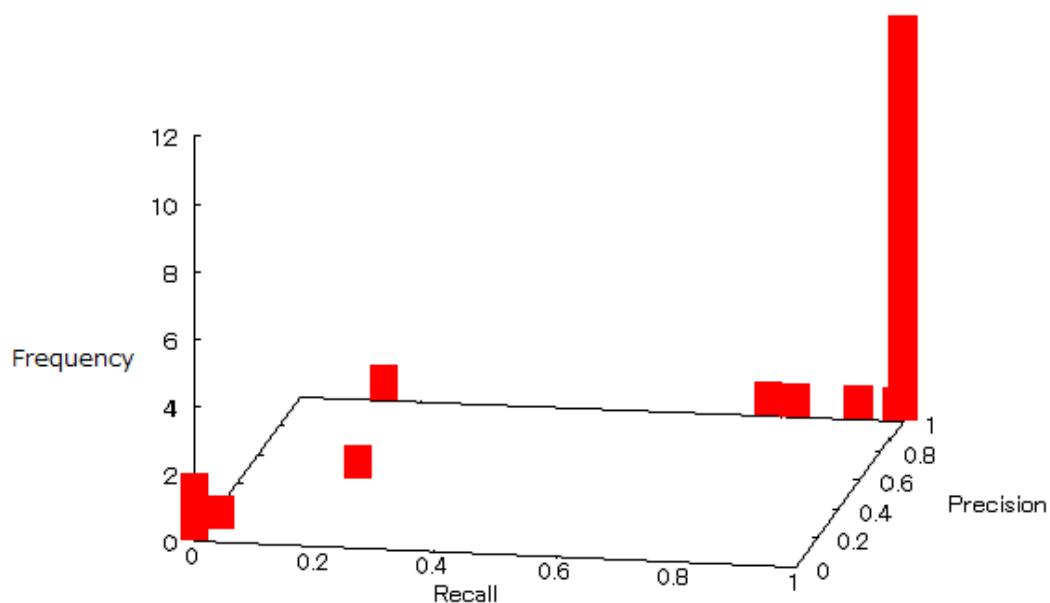


図 5.5 確定問合せの検索結果の再現率と精度

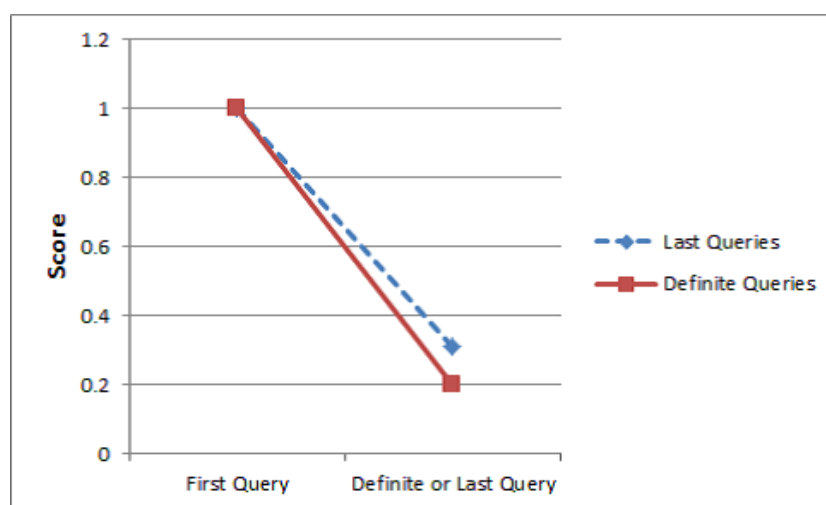


図 5.6 最初の間合せと確定問合せと最後の間合せのそれぞれのスコアの平均

間合せのスコア. 同じ間合せ記述プロセス内の間合せであっても、間合せ $q_{i,j,k}$ 同士を比較することは難しい。なぜなら、各間合せの検索結果は異なるスキーマ構造を持っており、再現率と精度を計算することが困難であるからである。従って、各 $q_{i,j,k}$ 毎に、その間合せがどれだけ正解間合せ cq_i と違っているかを表すスコア (difference score) を付けることを考えた。このスコアは関数 $ds_{cq}(q)$ で表す。このスコアは間合せ間の違いが小さいほど値が小さい。すなわち、 $q = cq$ の時、 $ds_{cq}(q) = 0$ であり、 q と cq の違いが大きければ大きいほど $ds_{cq}(q)$ は大きくなる。スコアの求め方の詳細は付録 B に示す。

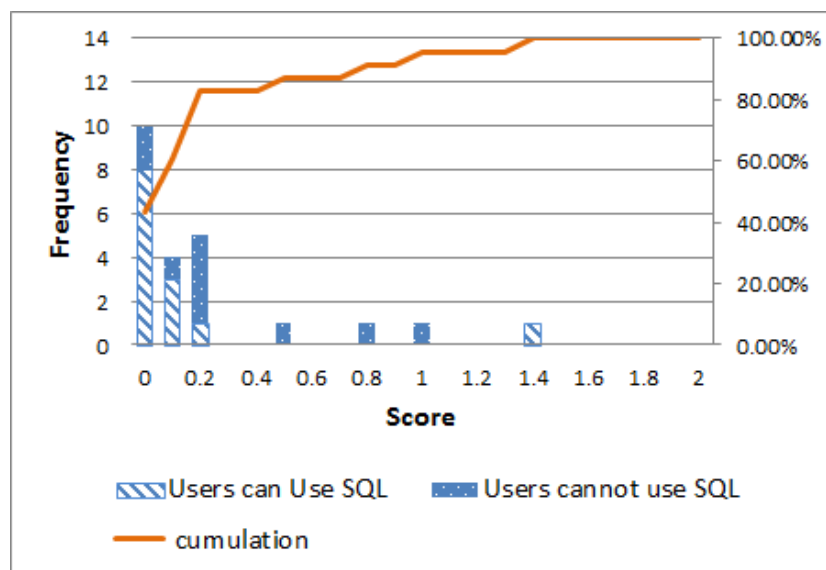


図 5.7 確定問合せのスコアの分布

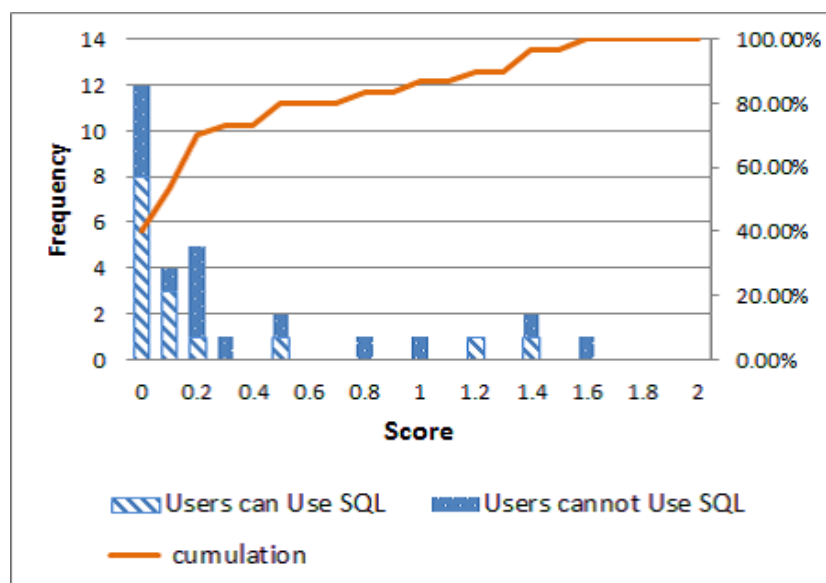


図 5.8 最後の問合せのスコアの分布

問合せ記述プロセス毎に、被験者が最初に行った問合せ、確定問合せ、最後の問合せのスコアを求めた。各スコアは最初に行った問合せのスコアが 1 になるように正規化した。各スコアの平均を図 5.6 に示す。23 個の確定問合せのスコアの平均は 0.199 であった。最後の問合せにおいても、スコアの平均は 0.310 であった。このことから、各被験者は最終的には正解問合せに近い問合せを記述できたことが分かる。

図 5.7 と 5.8 にスコアの詳細を示す。この図は、確定問合せ (図 5.7) と最後の問合せ (図 5.8) のス

コアのヒストグラムと累積度数分布である。確定問合せのうち 82% はスコアが 0.2 以下であった。また、SQL を使うことができない被験者によって書かれた確定問合せのうち 70% はスコアが 0.2 以下であった。更に確定問合せでないものも含む全ての最後の問合せに関しても、その 70% はスコアが 0.2 以下であった。まとめると、これらの実験結果は、平均的には、MorphingAssist はユーザがより精確な問合せ (正解問合せに近い問合せ) を記述することの支援に成功していることを示している。

図 5.8 では、最後の問合せに対するスコアを示したが、実際には、ユーザが書き直した問合せが適切でないと感じた場合は、ユーザは検索の際には元の問合せを使えばよく、最後の問合せを使う必要はない。

第 6 章

おわりに

本論文では、キーワード問合せと構造化問合せをシームレスに融合した問合せ言語を利用して、ユーザによる構造化問合せの記述を支援する手法 **MorphingAssist** を提案した。具体的には、キーワード問合せを入力したユーザに対して、そのキーワード問合せを構造化問合せに変更するためのヒントを提示する。ユーザは必ずしも構造化問合せが作成できるまで **MorphingAssist** を使う必要はなく、必要な範囲、出来る範囲で詳細化できればよく、この点において **MorphingAssist** は **pay-as-you-go** スタイルの構築支援を行っていると言える。

実データを用いた実験では、次の 2 つの観点から評価を行った。(1) 生成したヒントの質。(2) **MorphingAssist** の効果。実験では、**MorphingAssist** は適切なヒントを表示することに関してうまく働き、また **MorphingAssist** を使ったユーザはより精確な問合せが記述できた。

参考文献

- [1] DBpedia: wiki.dbpedia.org : About, <http://www.dbpedia.org/>, (参照 2013-01-15).
- [2] ACM DIGITAL LIBRARY, <http://dl.acm.org/>, (参照 2013-01-15).
- [3] dblp.rkbexplorer.com, <http://dblp.rkbexplorer.com/>, (参照 2013-01-15).
- [4] W3C: LinkingOpenData, <http://www.w3.org/wiki/SweoIG/TaskForces/CommunityProjects/LinkingOpenData>, (参照 2013-01-15).
- [5] 袖山広輝, 只石正輝, 安永ゆい, 品川徳秀, 森嶋厚行. “構造/テキスト Web データのためのハイブリッド問合せ言語”. Web とデータベースに関するフォーラム 2011, 9 pages, (2011).
- [6] Yui Yasunaga, Atsuyuki Morishima, Hiroki Sodeyama and Masateru Tadaishi. “Gradation: A Pay-as-You-Go Style Hybrid Query Language for Structured and Text Data”, Proceedings of iConference 2013, pp.209-220, (2013).
- [7] Demidova, E., Zhou, X. and Nejdl, W. “IQP: Incremental query construction, a probabilistic approach”, 26th IEEE International Conference on Data Engineering (ICDE 2010), pp.349-352, (2010).
- [8] Chaitanya Mishra and Nick Koudas. “Interactive Query Refinement”, Proceedings of the 12th International Conference on Extending Database Technology (EDBT 2009), pp.862-873, (2009).
- [9] internet.com K.K. (Japan) “An article from internet.com K.K.(Japan)”, <http://japan.internet.com/wmnews/20090908/5.html>(2009), (参照 2013-01-15).
- [10] Wikimedia Downloads, <http://dumps.wikimedia.org/>, (参照 2013-01-15).
- [11] DBpedia Japanese, <http://ja.dbpedia.org/>, (参照 2013-01-15).

本論文に関連する発表論文

国際会議論文

- Yui Yasunaga, Atsuyuki Morishima, Hiroki Sodeyama and Masateru Tadaishi. “Gradation: A Pay-as-You-Go Style Hybrid Query Language for Structured and Text Data”. Proceedings of iConference 2013, pp.209-220, Texas, Feb. 2013.

国内査読付き論文

- 袖山広輝, 只石正輝, 安永ゆい, 品川徳秀, 森嶋厚行. “構造/テキスト Web データのためのハイブリッド問合せ言語”. Web とデータベースに関するフォーラム (WebDB Forum)2011, 9 pages, 東京, 2011 年 11 月, (WebDB Forum 企業賞 (サイボウズ、サイボウズ・ラボ賞) 受賞).

国内研究会論文

- 安永ゆい, 森嶋厚行, 袖山広輝. “構造/テキスト Web データを対象とした Pay-as-You-Go スタイルの問合せ構築支援手法”. 第 5 回データ工学と情報マネジメントに関するフォーラム (DEIM2013). 福島, 2013 年 3 月.
- 安永ゆい, 袖山広輝, 森嶋厚行, 只石正輝. “構造/テキストデータ用ハイブリッド問合せ言語を用いた問合せ構築支援”. 情報処理学会第 74 回全国大会講演論文集 (第 1 分冊), pp.597-598, 名工大, 2012 年 3 月.

口頭発表

- Yui Yasunaga, Atsuyuki Morishima, Hiroki Sodeyama, Masateru Tadaishi. “Gradation: A Hybrid Query Language for Structured and Text Data”. Information Science Workshop2012, University of Tsukuba, Dec. 2012.

謝辞

本研究を進めるにあたって、本当に多くの方々に助けていただきました。

森嶋先生には、沢山のことを教えていただきました。もう間に合わない、そう思っていたのに先生に相談すると道が拓けたあの時 (IPSJ に向けての実験の時) は、本当に「物事には方法がある」んだ！と心から思いました。国際会議に投稿する前は本当に毎日沢山の時間をかけて論文を見て下さり、ありがとうございました。時には発破をかけてやる気を出させて下さいましたし、ただ単に研究を進めることだけではなく、様々な力をつけさせようと、勉強と成長の環境を用意して下さいたことにもとても感謝しております。もしもう一度大学3年生に戻っても、私はやはり森嶋研を選びたいです。

品川先生は、WenbDBf に投稿する前には、出張中であるにもかかわらず毎日多くの時間を添削に費やして下さいました。また普段から技術的な質問にも沢山答えて下さり、頼らせて下さいました。ありがとうございました。

杉本先生、阪口先生、永森先生、品川先生には、合同ゼミで様々な視点からのコメントをいただきました。ありがとうございました。ゼミ後のちょっとした時間にして下さるお話も興味深かったです。

研究室の先輩である只石先輩、袖山先輩には共同で研究を進めるにあたって、本当に助けていただきました。M1 の夏に共同で論文を執筆した怒濤のような日々は忘れることができません。WebDBf で企業賞ももらうことができ、嬉しかったですね。只石先輩は、休学なさってからも、私の質問に丁寧に答えて下さったり、研究にアドバイスを下さったりと、大変お世話になりました。袖山先輩には、お世話になりっぱなしでした。M1 の頃は研究打合せにもほぼ毎回参加して下さい、IPSJ74 や DEIM2013 の論文執筆の際は多くのコメントを下さりました。そして袖山先輩の研究なしには私の研究も存在しませんでした。本当にありがとうございました。

卒研の頃から含めて、実験に協力してくれたサークルの仲間、学類・研究科の皆様にもとても感謝しています。毎回急なお願いになってしまったのにもかかわらず快く引き受けてくれた皆のおかげでどの実験も無事期間内に行い、研究をまとめることができました。

弓矢先輩、太田先輩、山潤先輩、「先輩が居る」あの頃がとても懐かしいです。1年上の先輩方には、研究室のことやオセロプログラミング、珠玉ゼミ、自主ゼミ、卒研、修士に入ってからの研究と、本当に沢山のことを教えていただきました。先輩方の優しいアドバイスと的確なコメントに、私

は何度もほっとし助けていただきました。

同期である三津石君とは、森嶋研に配属されてから約3年、励まし合ったり、一緒に反省したりと、とてもよい仲間になりました。Datalogのプレゼンをしたあの頃よりはお互い成長したなあと思うよね。同期が三津石君で良かった！

研究生生活を応援してくれた家族にもとても感謝しています。なかなか連絡をしない私でたくさん心配をかけたと思いますが、やりたいと言ったことを応援してくれた父、明るいメールを沢山くれた母、すてきな小物を沢山作ってくれた妹に感謝しています。お米やみかんを育てて送ってくれる父方のおじいちゃんおばあちゃん、ありがとうございました。食べ物から沢山元気をもらいました。遊びに帰る度にリラックスと元気の充電をさせてくれる母方のおばあちゃんと叔母にも、心からありがとう。

秘書である篠崎さん、出張の手続きだけでなく月に1度のアメニティ大会もありありがとうございました。お家でかき氷を作ったりたこ焼きを焼いたり、そういうお母さんになりたいなあと思いました。多和田さんの手作りお菓子もとてもおいしかったです。多和田さんのような素敵なお母さんもあこがれます。

朝来ると、ひとこと話しかけてくれるお掃除のおばちゃんもありありがとうございました。ラウンジも廊下も流しもいつもきれいで、毎日気持ちよく過ごせました。顔を覚えていてくれて、行くと必ず笑って話しかけてくれる食堂のおばちゃんのおかげで、ちょっと元気がない日も、笑顔になれました。

研究室のメンバーとして楽しい時間と研究への沢山の指摘、それから実験結果解析のお手伝いをしてくれた後輩のみんなにも感謝しています。青木君、コツコツと積み上げていく青木君の姿勢は見習いたいなあと常々思っていました。福角君、さっと仕事を片付けて、いつも机の上がキレイで、研究も遊びも楽しんでいる姿は私の理想です。愛ちゃん、色んな話題を振って、研究室を明るくしてくれてありがとう。私にない強い心を持っていて、すごいなあと思っています。富田さん、卒研を始めた頃は研究室内に詳しい人がいなくて大変そうだったのに、最後にはきちんと実験結果を出していて、すごいなあ。卒業してしまう前にもっとたくさんお話してみたいなあと思っています。権ちゃん、権ちゃん先生にはとてもとてもお世話になりました！先生がいなくてはきっとシステムも完成しませんでした。丹治君、バリバリっとシステムを作ってちゃんと実験も予定通りにやって、ほんとにすごいなあ。実は丹治君の受け答えに癒されていました。

望月先輩は、研究が進んで忙しくしている時も、行き詰まって落ち込んでいるときも、いつも暖かく見守って励ましてくれました。沢山話を聞いてくれて、アドバイスをしてくれて本当にありがとうございました。

本当に、皆様のおかげで、この論文を書くことができたのだなあと心から思います。ありがとうございました。

付録 A

実験で用いた情報要求と正解問合せ

実験で使用した情報要求とそれと対応する正解 Gradation 問合せのリストは表 A.1 の通りである.

表 A.1 情報要求とそれと対応する正解 Gradation 問合せのリスト

i	情報要求 n_i 正解 Gradation 問合せ cq_i
1	神奈川にある公園を全て知りたい <code>dcterms:subject == "category-ja:神奈川県のパーク"</code>
2	東京都内にある医療機関全てと, その医療機関の所在地が知りたい <code>(dcterms:subject == "category-ja:東京都の医療機関") [t_uri, dbpprop-ja:所在地]</code>
3	ロンドンがホームタウンであるサッカークラブとそのクラブのホームスタジアムが知りたい <code>(dbpprop-ja:ホームタウン == "dbpedia-ja:ロンドン") [t_uri, dbpprop-ja:スタジアム]</code>
4	日本にある灯台を全て知りたい <code>dcterms:subject == "category-ja:日本の灯台"</code>
5	日本にある女子高を全て知りたい <code>dcterms:subject == "category-ja:日本の女子高等学校"</code>
6	日本のロックバンドと, そのバンドが所属しているレコード会社が知りたい <code>(dcterms:subject == "category-ja:日本のロック・バンド") . dbpprop-ja:label . (dcterms:subject == "category-ja:日本のレコード会社")</code>
7	1990 年生まれのタレントを全員知りたい <code>dcterms:subject == "category-ja:日本のタレント" dbpprop-ja:生年 == 1990</code>
8	日本にある保守政党全てと, その政党の成立年が知りたい <code>(dcterms:subject == "category-ja:日本の保守政党") [t_uri, dbpprop-ja:成立年月日]</code>
9	つくば市にある研究所を全て知りたい <code>dcterms:subject == "category-ja:つくば市の研究所"</code>
10	イタリアセリエ A に所属するサッカークラブ全てと, その監督名が知りたい <code>(dbpprop-ja:リーグ == "dbpedia-ja:セリエ A_(サッカー)") [t_uri, dbpprop-ja:監督]</code>

付録 B

スコア $ds_{cq}(q)$ の求め方

本付録では、ユーザが書いた問合せ q と正解問合せ cq が与えられた時に、 q と cq の違いを表すスコア (difference score) を計算するための関数 $ds_{cq}(q)$ について説明する。この関数は 5.4.2 節の評価で利用したものであり、 q と cq の違いが小さいほど値は小さくなる。すなわち、 $q = cq$ の時、 $ds_{cq}(q) = 0$ であり、 q と cq の違いが大きければ大きいほど $ds_{cq}(q)$ は大きくなる。

スコア $ds_{cq}(q)$ は 図 B.1 で求められる。具体的には、まず q 中の問合せ構成要素 $c_i \in component(q)$ に対して、それと対応する正解 cq 中の構成要素 $cc_j \in component(cq)$ を見つけ、その 2 つがどれくらい異なるかによってペナルティ (c-penalty) を加点する。次に、 cq 中の構成要素 $cc_j \in component(cq)$ のうち、先の作業で使われなかった要素の数だけペナルティ (m-penalty) を加点する。図 B.1 の 4-8 行目で、c-penalty を加点している。ここで、5 行目の $counterpart_{q,cq}(c_i)$ は、 c_i に対応する cq 中の構成要素を求める関数である。この詳細は B.1 節で述べる。また、6 行目の $c-penalty(c_i, cc_j)$ は構成要素の差異に応じた c-penalty を求める関数である。c-penalty の点数については B.2 節で述べる。2,7,10 行目で m-penalty を加点している。10 行目の $m-penalty$ は定数であり、今回は 5 としているが、この点数についても B.2 節で述べる。

B.1 構成要素の対応付け関数

$counterpart_{q,cq}$ は、ユーザが記述した問合せ q と 正解問合せ cq のペアがある時に、 q 中の構成要素 c_i が、 cq 中のどの構成要素の書き換えであるか (どの構成要素を意図して書いたか) を表す関数である。例えば、情報要求が「トムという俳優を全員知りたい」だった場合に、 $q = \text{"tom actor"}$, $cq = \text{"name==tom type==actor"}$, $c_1 = \text{"tom"}$, $cc_1 = \text{"name==tom"}$ がある時、 $counterpart_{q,cq}(c_1) = cc_1$ が成り立つ。関数 $counterpart_{q,cq}$ は、 q の構成要素の集合 $component(q)$ を定義域とし、 cq の構成要素 $component(cq)$ と nil からなる集合を値域とする。一般には、問合せ作成者の意図をくんで $counterpart_{q,cq}$ を計算するのは困難であるが、上記のように定義域と値域を限定する事により、「 cc_j が c_i の類義語を含む」時に $counterpart_{q,cq}(c_i) = cc_j$ が成立し、そうでない場合は $counterpart_{q,cq}(c_i) = nil$ が成立する、というヒューリスティクスが利用可能な場合が

INPUT: query q , correct query cq

OUTPUT: difference score $score$

```

1:  $score = 0$ ;
2:  $rest = component(cq)$ ;
3: // penalty for each component を計算
4: FOREACH  $c_i \in component(q)$ 
5:    $cp_i = counterpart_{q,cq}(c_i)$ 
6:    $score = score + c-penalty(c_i, cp_i)$ ;
7:  $rest = rest - cp_i$ 
8: ENDFOR
9: // penalty for missing components を計算
10:  $score = score + m-penalty * |rest|$ 

```

図 B.1 $ds_{cq}(q)$ の計算

ある．今回の実験で利用した全ての問合せは，このヒューリスティクスが適用可能なケースであったため，これを利用して $counterpart_{q,cq}$ を計算した．

B.2 ペナルティの値の決定

$c-penalty(c_i, cp_i)$ および $m-penalty$ の値は次の原則に従い決定した．

原則： c_i の存在によって q の検索結果が cq の検索結果により近くなる場合にペナルティの値はより小さくなるべきである．「検索結果に近い」ということを定めるために次の 3 つのヒューリスティクスを用いる．

(ヒューリスティック 1) c_i がキーワードであるとき，(a) cp_i がその一部に文字列 c_i を含む場合，(b) cp_i がその一部に文字列 c_i の類義語を含む場合，(c) cp_i が nil である場合，の順に q の検索結果は cq の検索結果に近くなる．例えば， q の構成要素 c_i が *thomas* であるときよりも，*tom* であるときの方が， q の検索結果は $name==tom$ を含む正解問合せ cq の検索結果に近くなると考えられる．

(ヒューリスティック 2) c_i がキーワード以外の構成要素である場合は，(a) $result(c_i) = result(cp_i)$ ，(b) $result(c_i) \supset result(cp_i)$ ，(c) $result(c_i) \subset result(cp_i)$ ，(d) $cp_i = nil$ の順に q の検索結果は cq の検索結果に近くなる（ここで $result(c_i)$ は c_i 単体の検索結果を表す）．例えば， $type==japaneseActor$ よりも $type==allActor$ が問合せに含まれる方が， $type==asiaActor$ を含む問合せの検索結果に近くなると考えられる．

(ヒューリスティック 3) cp_i と同じ構成要素が $component(q)$ 中に存在する場合の方が，そうでない場合よりも q の検索結果は cq の検索結果に近くなる．例えば，次の 2 つの場合の検索結果を考える．まず， $q = \text{“tom actor”}$ ， $c_1 = \text{“tom”}$ ， $cq = \text{“name==tom type==actor”}$ ， $cc_1 = \text{“name==tom”}$

表 B.1 c-penalty の値

type of c_i	relation between c_i and cp_i	$cp_i \in \text{component}(q)$	$cp_i \notin \text{component}(q)$
keyword	$match(c_i, cp_i)$	1	3
	$match(similar(c_i), cp_i)$	2	4
	$otherwise (cp_i = nil)$	5	5
non-keyword component	$result(c_i) = result(cp_i)$	0	0
	$result(c_i) \supset result(cp_i)$	1	2
	$result(c_i) \subset result(cp_i)$	3	3
	$otherwise (cp_i = nil)$	5	5

であるとする．この時， $counterpart_{q,cq}(c_1) = cc_1 = \text{“name==tom”}$ は q 中に存在しない．次に， $q' = \text{“tom name==tom actor”}$ ， $c'_1 = \text{“tom”}$ であるとする．この時， $counterpart_{q',cq}(c'_1) = cc_1 = \text{“name==tom”}$ は q' 中に存在する． q と q' の検索結果を比べると， q' の方が cq の検索結果に近くなると考えられる．

上記の原則およびヒューリスティクスに矛盾しないように考慮した結果，c-penalty を表 B.1 とし，m-penalty=5 とした．表 B.1 の行の場合分けにヒューリスティクス 1 と 2 を使い，列の場合分けにヒューリスティック 3 を用いている．